

WEEK 1:

~~Hyperparameter Tuning, Batch Normalisation,
Programming frameworks~~

Practical aspects of Deep Learning

Applied ML is highly iterative process

- # layers

- # hidden units

Learning rates

activation functions

NLP, CV, Speech recog., Structured Data

→ Train / dev / Test sets

Data

Training set

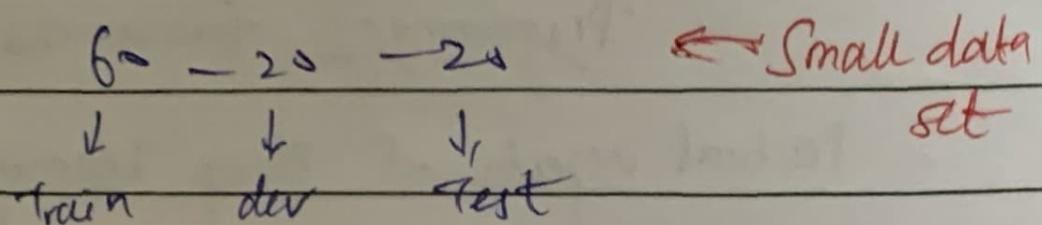
- Hold out validation
- Dev set

Test

Train

Date:

Previous era: 70 → 30 ← Test



Big data era:

[dataset just chooses best
choices of different algorithms, which is
working best]

98 - 1 - 1

← Big dataset

99.5 - 0.25 - 0.25

99.5 - 0.4 - 0.1

Date:

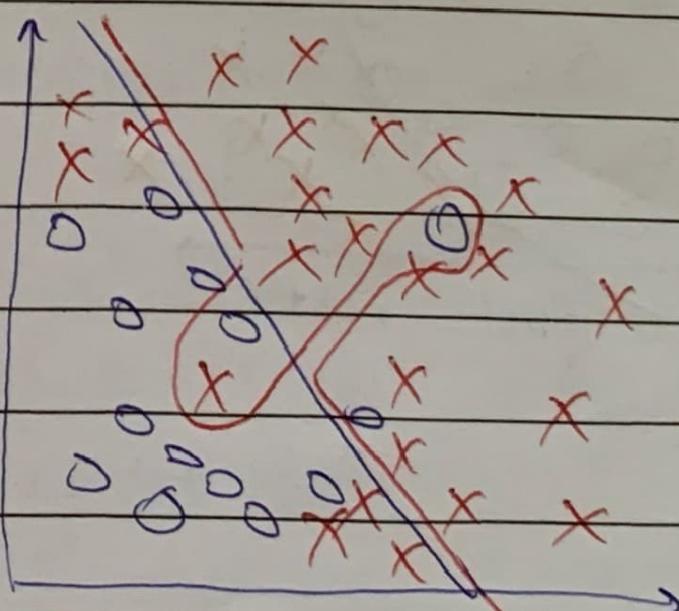
→ Mismatched train/test distribution :

Make sure the dev and test set come from the same distribution.

→ Not having a test set might break.

∴ Test set gives you the unbiased estimate of how your final network works.

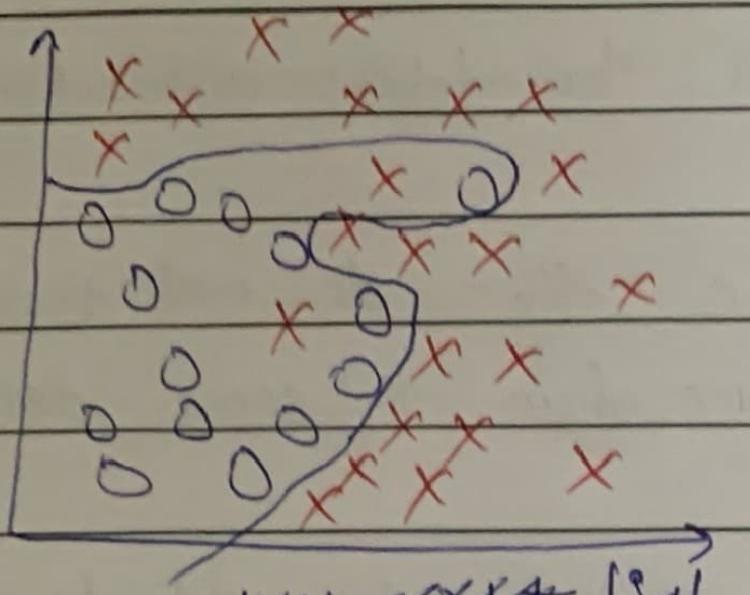
→ Bias and Variance



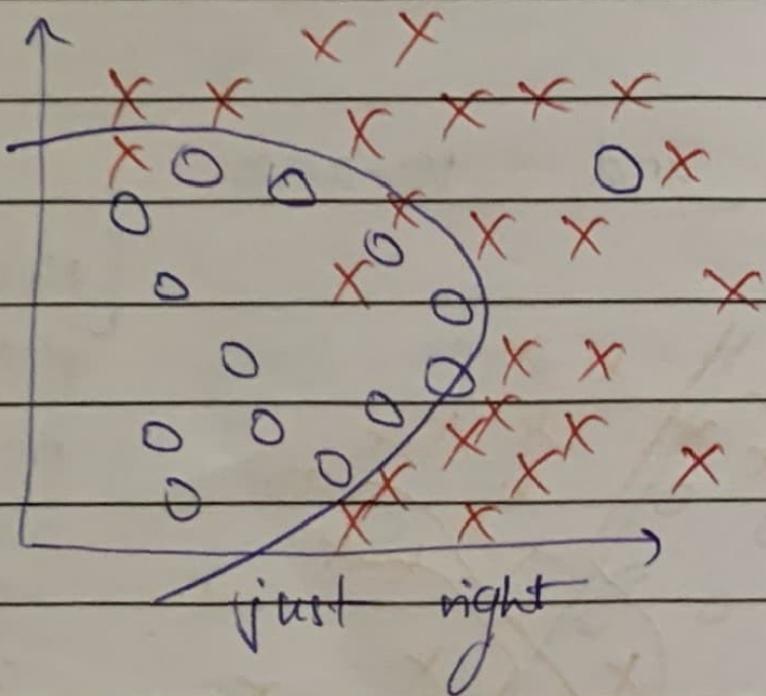
{ Note: The red plot boundary is high bias & high variance
∴ it is mostly linear and middle two mislabelled are covered }

high bias
"underfitting"

Date:



"overfitting"



Date:

Example:
→ Cat classification:

Train set error: 1% | 5%

Dev set error: 11% | 16%

high variance | high bias

Human ≈ 0%.

Optimal (Bayes) error ≈ 0%. | 15% | 0.5%

30% | 1%

high bias | low bias
high variance | low variance

The analysis depends on the optimal error.

→ Basic "recipe" for ML

High Bias?

(Training data perform.)

↓ N

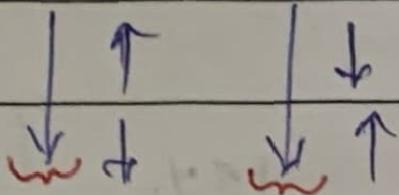
- ① Bigger network
- ② Train longer
- ③ (NN architecture search)

Date: _____

High variance? → ① more data
(dataset perform.) ↘ ② Regularisation
↓ ③ NN architecture selection

{ Done }

"Bias Variance trade-off"



Deep Learning to decrease both

"Training a bigger NN almost never hurts."

Date:

⇒ Regularization (helps reducing high variance)

logistic Regression

$$\min_{w,b} J(w,b)$$

$w \in \mathbb{R}^n$; $b \in \mathbb{R}$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$+ \underbrace{\lambda}_{\text{regularization parameter (hyperparameter)}} \frac{\|w\|_2^2}{2m}$$

L_2 regularization (most common type):

$$\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$$

{Note: $+ \frac{\lambda}{2m} b^2$ can also be done, but not used}

∴ b is just a number
 w has all the features in it}

Note: we will write "lambda" instead of "lambda". Date: "lambda".

L regularization

$$\frac{1}{m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

Here, w will be sparse.

NN:

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)})$$

$$= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l-1)}} (w_{ij}^{(l)})^2$$

$$w = \begin{pmatrix} n^{(1)} \\ \vdots \\ n^{(L-1)} \end{pmatrix}$$

"frobenius norm"

Date:

$$dw \rightarrow (\text{from back prop.}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

(can be called as :

"weight decay" : $w^{[l]} * \left(1 - \frac{\lambda\alpha}{m}\right) \approx < 1$

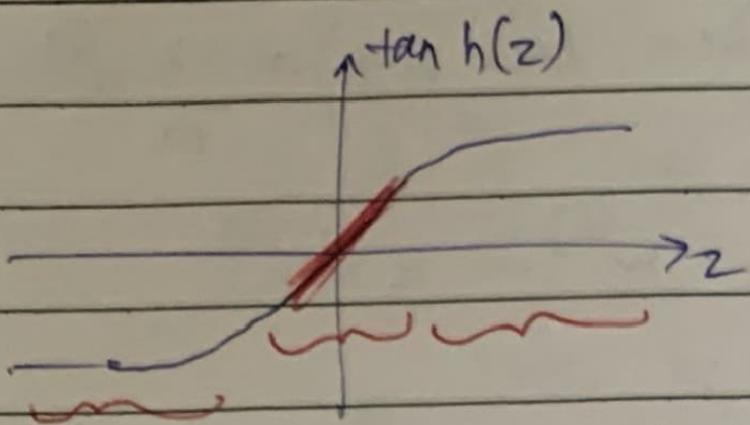
$$w^{[l]} := w^{[l]} - \alpha \left[(\text{from back prop}) + \frac{\lambda}{m} w^{[l]} \right]$$

$$\therefore w^{[l]} = w^{[l]} - \alpha \frac{\lambda}{m} w^{[l]}$$

$-\alpha$ (from backprop)

If the value of lambda is set to "very high"
then the weight matrix $w^{[l]}$ is zeroed
and high variance turns to high bias.
this causes zeroing or reducing impact
of hidden units. Many of The hidden units don't make
an impact and hence the ^{Deep}NN converts to a
smaller NN.

Date: 26/07/2024



$$\lambda \uparrow \quad w^{[l]} \downarrow \quad z^{[l]} : \quad w^{[l]} a^{[l-1]} + b^{[l]}$$

$z^{[l]} \downarrow \quad \vdash g(z^{[l]})$ will roughly
linear since the $\tanh(z)$
at small z values is almost linear.

Date:

→ Dropout Regularisation:

For each training example, the hidden units are diminished with a given probability and then it is trained.

Implementing Dropout ("Inverted Dropout")

Illustrate in layer $l=3$:

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0],$
 $a_3.\text{shape}[1])$
< keep-prob

keep-prob = 0.8

0.2 chance of
eliminating any
hidden unit.

$a_3 = \text{np.multiply}(a_3, d_3)$

$\# a_3 = d_3$

Date: _____

$$a_3' = \text{keep-prob}$$

Note: This step ensures final value of a_3' is still same so that ~~there~~ there is no scaling problem.

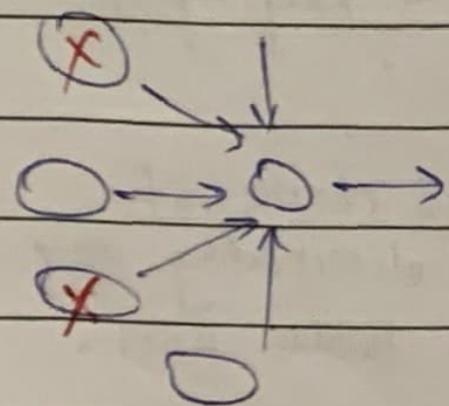
Making predictions at test time:

No drop-out.

why does drop-out work?

→ Can't rely on any one feature, so have to spread weights.

→ Shrink weights



Drop-out and L_2 regularisation have almost the same effect.

Date: _____

You could keep the value of keep-prob different for different layer, depending on the size of the weight matrix.

For a ~~higher~~ big weight matrix; keep-prob could be 0.5 otherwise keep-prob can be probably higher than 0.5.

The calculation of cost function $J(w, b)$ becomes difficult along with the drop out.

→ Regularization methods:

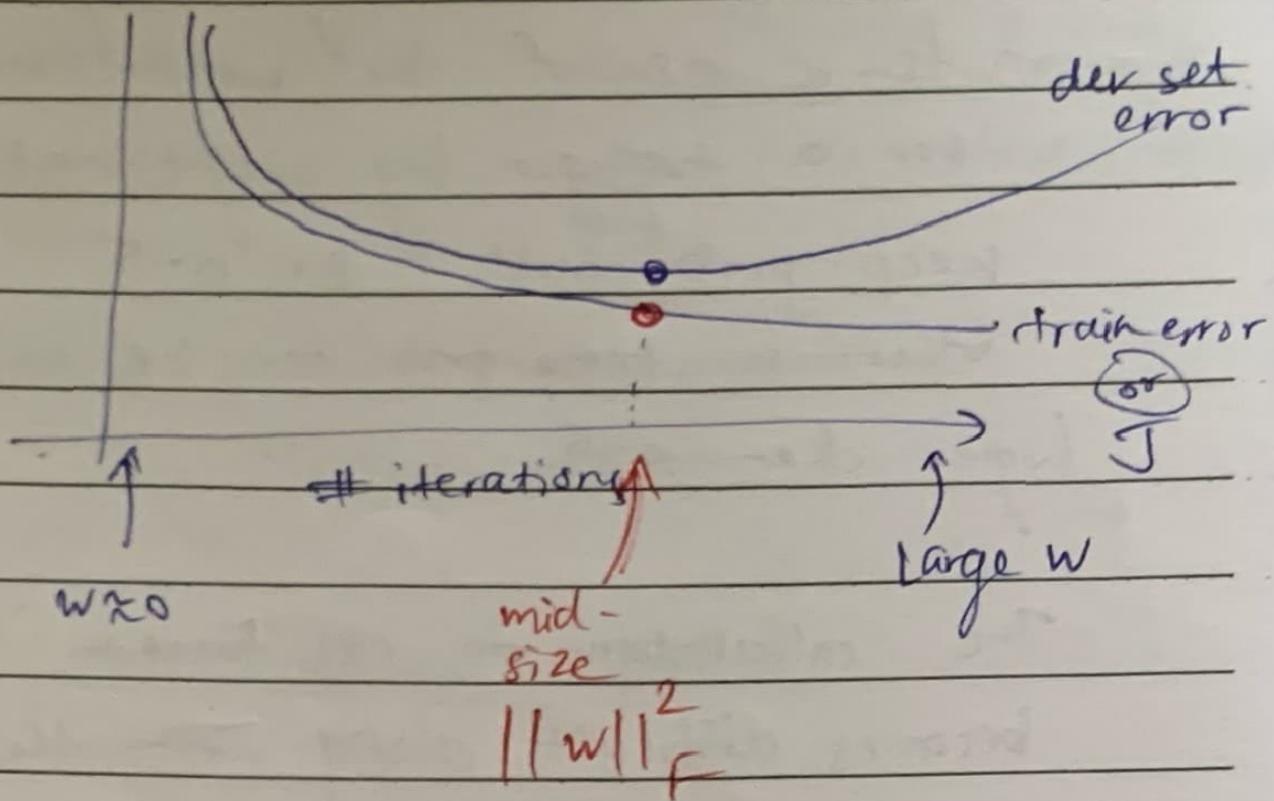
→ Data augmentation:

Example: To prevent overfitting, in a cat image classification, flip the images horizontally, which will get 2m training examples.

Thus, can get more data by augmenting the images.

Date:

→ Early Stopping (stops the training
~~starting early~~)



Early stopping carries out
the processes together
called as orthogonalisation.

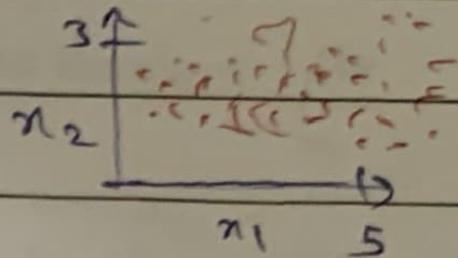
- ① Optimise cost fn J
- ② Not overfit

hence to stop early it will
even stop the minimising of
cost function J.

Date:

Normalizing training sets :

$$X = \begin{bmatrix} n_1 \\ n_2 \end{bmatrix}$$

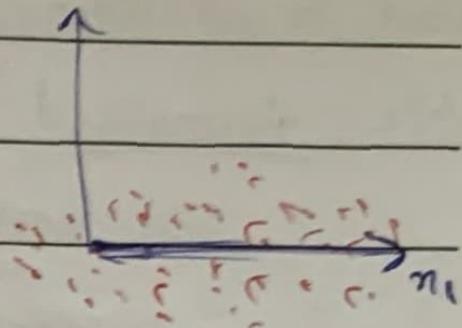


Steps:

①

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$X := X - \mu$$



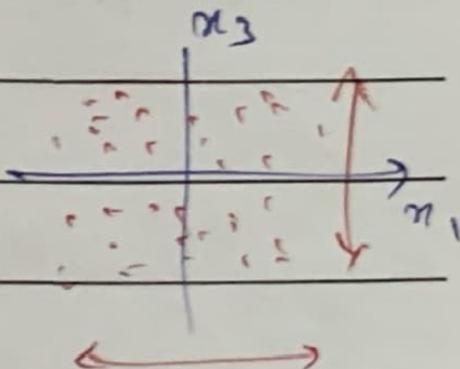
② Normalize variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \# 2$$

vector of variances ↓ element-wise

~~vector~~ of each of the features

$$X / = \sigma$$



Date: _____

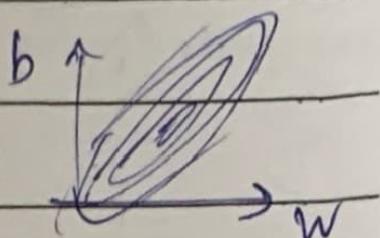
Use same μ , σ to normalise test set.

why normalize inputs?

while Unnormalised: The cost function is

very squished ~~and~~

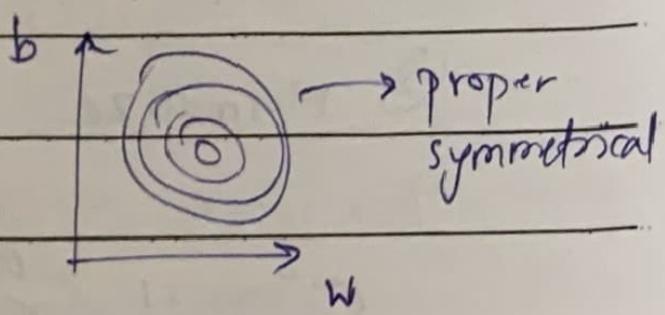
Need to use ^{small} learning rate



while normalised:

cost function on average will look
symmetric.

can be use ↑ learning
rate.



→ Easier to optimise when normalised.

Date:

Vanishing / Exploding gradients:

$$x_1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

$$x_2 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad \rightarrow y$$

$$w^{[1]} \quad w^{[2]} \quad \dots \quad w^{[L]}$$

$$g(z) = z, \quad b^{[l]} = 0$$

$$\hat{y} = w^{[L]} \cdot w^{[L-1]} \cdot w^{[L-2]} \cdots w^{[2]} w^{[1]} x$$

$$z^{[1]} = w^{[1]} x$$

$$a^{[1]} = g(z^{[1]}) = z^{[1]}$$

$$a^{[2]} = g(z^{[2]}) = g(w^{[2]} a^{[1]})$$

⋮

$$w^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$\hat{y} = w^{[l]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} x$$

~~$$\hat{y} = w^{[l]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} x$$~~

~~$$\hat{y} = w^{[l]} 1.5^{l-1} x$$~~

Date: _____

∴ If \rightarrow Deep NN very deep

L is very large

then the gradients will explode;

then the γ will be very large

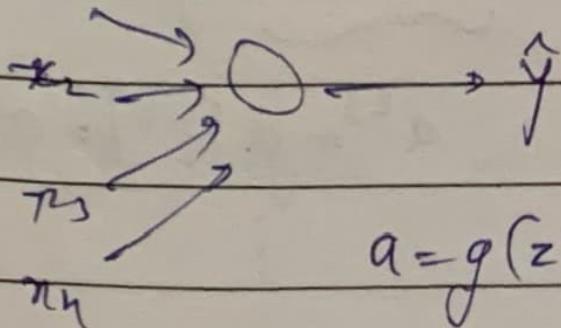
$w^{[L]} > I$ (activations explode)

$w^{[L]} < I$ (activations decrease [vanishing exponentially])

→ Weight initialisation for DNN:

$a^{[1]}$

x_1



$$a = g(z)$$

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

large $n \rightarrow$ small w_i

Date:

$$\text{var}(w_i) = \frac{1}{n}$$

$$w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt} \left(\frac{1}{n^{[l-1]}} \right)$$

$$\text{If } g^{[l]}(z) = \text{ReLU}(z)$$

$$\text{then } \text{var}(w_i) = \frac{2}{n}$$

$$\text{and } \text{np.sqrt} \left(\frac{2}{n^{[l-1]}} \right)$$

Other variants:

$$\text{for } \tanh(): \text{np.sqrt} \left(\frac{1}{n^{[l-1]}} \right)$$

Xavier Initialisation

other version

$$\text{np.sqrt} \left(\frac{2}{n^{[l-1]} + n^{[l]}} \right)$$

Date: _____

Gradient check for a NN:

Take $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]} \dots$

and reshape into

big vector θ

$$J = (w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

Take $d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$

and reshape into a

big vector $d\theta$

Gradient checking (road check):

$$J(\theta) = J(\theta_1, \theta_2, \dots)$$

for each i :

$$d\theta_{approx}^{[i]} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon)}{2\epsilon}$$

[i]

Date:

$$\approx d\theta \frac{\partial \theta}{\partial i} = \frac{\partial J}{\partial \theta_i}$$

$$d\theta_{\text{approx}} \approx d\theta$$

Check Euclidean distance

$$\frac{\| d\theta_{\text{approx}} - d\theta \|_2}{\| d\theta_{\text{approx}} \|_2 + \| d\theta \|_2}$$

Euclidean lengths ← This is calculating

Euclidean distance

formula : $\sqrt{\sum_{i=1}^n (d\theta_{\text{approx}} - d\theta)^2}$

$$\epsilon = 10^{-7}$$

if ans = 10^{-7} → gradient ✓
→ great!

if ans = 10^{-5} → double check

if ans = 10^{-3} → worry (bug)

$> 10^{-3}$ → worry

derivatives → might be wrong

Date:

→ Gradient checking Implementation Notes:

① Don't use while training;

use only for debugging

② If algorithm fails grad check;

look at components to try
to identify bug. Check for $(\mathbf{d}b^{[l]}, \mathbf{dw}^{[l]})$

③ Remember regularization.

④ Doesn't work with dropout

⑤ Run at random initialisation, to check
if the it work values of $w^{[l]}, b^{[l]}$
 $\neq 0$.